# Comparative Study of Profiling Tools on Fugaku Supercomputer

Samar A. Aseeri<sup>1</sup>, Judit Gimenez<sup>2</sup>, Sameer S. Shende<sup>3</sup>, Benson K. Muite<sup>4</sup> and David E. Keyes<sup>1</sup>

1) KAUST Extreme Research Computing Center (ECRC), Saudi Arabia

2) Barcelona Supercomputing Center, Technical University of Catalonia, Spain

3) University of Oregon, USA 4) Kichakato Kizito, Kenya

February 26, 2023

This paper examines the overhead of performance tool measurements when profiling Fast Fourier Transform FFT based solvers for the Klein Gordon equation.

- This study is evaluating the performance of FFT algorithms on the Fugaku supercomputer, an application with many scientific use cases.
- Performance tools are being used to measure the improvements that a specific computer architecture can bring. This allows scientists and engineers to better understand how their programs will run on the given computer.

- Engineers are working to create bigger, advanced computers with specialized hardware known as accelerators. These accelerators can process tasks faster than data can be sent between processes. Therefore, in many cases, simulations must be run for long periods of time to solve moderate-sized problems.
- Scientific programs such as computational chemistry and fluid mechanics use libraries of basic components from the community. New research suggests that Dragonfly networks, a type of computer technology, could have a negative effect on FFT codes, which require a lot of synchronous communication between various components.

- Testing on the Fugaku supercomputer showed that performance of programs using FFT libraries was improved by minimizing interference to communication.
- Profiling and tracing techniques, such as binary instrumentation, were used to measure and analyze program performance. This requires knowledge of the program's structure and probing of memory set aside by the Linux kernel.
- Profiling using tracing has a higher overhead than traditional sampling, but the temporal information it reveals can be extremely helpful when analyzing programs without access to their source code.

The current study focuses on performance measurement using four tools and two FFT libraries, but there are a number of other tools and FFT libraries that can also be used. The profiled codes could be tuned to perform more efficiently on the given architecture. Performance tool overheads are minimized by selecting only a few measurements that may not be the most informative for other applications.

- 2 Background & Related Work
- 3 Profiling Results: A Comprehensive Overview
- 4 Evaluating the Results of Profiling
- 5 Summary & Conclusions



# Outline

#### Experimental Setup

- 2 Background & Related Work
- 3 Profiling Results: A Comprehensive Overview
- Evaluating the Results of Profiling
- Summary & Conclusions
- 6 Further Work

The Klein Gordon solver is used as a mini application benchmark for a parallel program that solves a partial differential equation. The FFT algorithm is one of the "Top 10 algorithms of the 20th century" and is simpler to use than alternative algorithms for solving linear equations. To ensure scalability, a two-dimensional pencil decomposition is used by the 2Decomp&FFT and FFTE libraries for the three-dimensional FFT algorithm. The grid size of 768<sup>3</sup> was chosen to allow examination of strong scaling behavior and does not utilize a huge amount of computational time.

Fugaku is an ARM architecture-based computer at the Riken Center for Computational Science in Kobe, Japan. It is one of the most powerful computers in the world and has been in operation since 2020. The Fugaku supercomputer consists of 158,976 homogeneous nodes consisting of Fujitsu A64FX CPU (48+4 cores) per node and has a performance of 442 PFLOPS measured by the LINPACK benchmark.

- TAU is an open-source set of tools that supports FORTRAN, C++, C, Java, Python, and UPC applications for performance measurement.
- Extrae is an open source tracing tool that can profile MPI, OpenMP, OMPss, pthreads, CUDA, OpenCL, OpenACC, GASPI and Python multiprocessing programs.
- FIPP and FAPP are Fujitsu vendor-provided tools used in this study to profile overall performance based on sampling.

- 2Decomp&FFT is a FORTRAN software library for parallel applications that utilize the Fast Fourier Transform or need to do neighbor data exchange communications on structured grids. It uses two MPI\_ALLTOALLV data exchanges to perform either a forward or backward parallel FFT. The memory requirements of the Klein Gordon solver using this library are dominated by 9 double complex arrays, 4 double complex arrays of combined size Nx × Ny × Nz and a temporary double complex array.
- FFTE is a FORTRAN only parallel Fast Fourier Transform library that uses Spiral to generate FFT kernels in C. It uses one MPI\_ALLTOALL exchange for a forward or backward parallel FFT. Memory requirements are dominated by 11 double complex arrays of size  $Nx \times Ny \times Nz/Np$  allocated in the main program.

- 2 Background & Related Work
- Operation of the second sec
- Evaluating the Results of Profiling
- Summary & Conclusions
- 6 Further Work

# Related Work

- Leu, et al. measured the overhead of performance tools Craypat, FPMPI, mpiP, Scalasca, and TAU on a Cray XC40 platform. The study found that FFTE had better strong scaling behavior than 2Decomp&FFT because it used fewer collective all to all communications.
- Cankur, et al. presents a comparative evaluation of several profiling tools Caliper, HPCToolkit, TAU and Score-P. Its results shows which tool produces the lowest overheads and most meaningful gathered information.
- Hunold et al. conducted a study comparing the overheads of six performance analysis tools on a single node and 32-node Intel Xeon Hydra cluster.
- A PRACE report compared the tools HPCToolkit, Extrae, SCALASCA and the Intel Trace Analyzer and Collector (ITAC).
- This work differs from other works in several ways, such as experimenting on an ARM architecture machine with a torus network.

- 2 Background & Related Work
- 3 Profiling Results: A Comprehensive Overview
- Evaluating the Results of Profiling
- Summary & Conclusions



- The FFTE library with Spiral generated C kernels with SVE instructions was the most performant in the current study.
- The Fugaku supercomputer offers a communication mode, the Small-Torus mode, which increases job queuing time and reduces job completion variation time. This mode is used for further experiments reported here.

# **Profiling Results**



Figure: Test Case Compared with TAU in the bottom plot in Small-Torus mode and without Small-Torus mode in the top plot

- TAU is used for sample based profiling of the compiled executable.
- TAU allows a user to observe events inside system libraries and we can see that the routine utofu\_read\_tcq\_nocb in the libtofucom.so library accounted for 3.37 seconds in this run.
- In this manner, TAU can show the precise time spent in MPI calls for each rank, the time spent in external libraries (like FFTW), the time spent inside MPI calls in proprietary libraries - all without modifying the application binary on Fugaku, by simply invoking the application using tau\_exec -ebs.

# TAU Profiling

| TAU: ParaProf: Statistics for: node 0 - RUN-2decomp/fugaku.ppk  |        |         |       |      |
|---|--------|---------|-------|------|
| Name  | Excl   | Incl. 7 | Calle | Chil |
|   | 19 425 | 27.69   |       | 252  |
| ▼ [CONTEXT] TAIL application  | 10.433 | 17 522  | 570   | 2.52 |
| CONTRACT, THE application   | 12 767 | 12 767  | 421   | 0    |
| ESIMMARY float 28 cost (U)(0)002 (hp210182 (data (handmark2)/version2decomp (K)ainCordon/2decomp tau (KoSamilmo2d 900)                  | 2 52   | 2 52    | 94    | 0    |
| SAMPL float 28 const [[/volood3/npz10139]/data/benchmark/?/volood3/npz20comp/klenGordon/zdecomp_tarkgemmips.rsoj]                       | 1.96   | 1.96    | 65    | 0    |
| [SAMP] = float128 const {//vol0003/hb210193/data/benchmark3/version2decomp/KleinGordon/2decomp tau/KaSemilma3d.f90} {223}               | 0.54   | 0.54    | 18    | ő    |
| [SAMP1] float 28 cost [[/vol003/hp210193/data/benchmark3/version2decom/KleinGordon/2decomp tai/kGsemims/dat9] [266]]                    | 0.02   | 0.02    | 1     | Ő    |
| [SUMMARY] enercalc [[/vol0003/hp210193/data/benchmark3/version2decomp/KleinGordon/2decomp_tau/enercalc.[90]]                            | 1.04   | 1.04    | 35    | 0    |
| SUMMARY storeold [//vol0003/hp210193/data/benchmark3/version2decomp/KleinGordon/2decomp_tau/storeold.f90]                               | 0.38   | 0.38    | 12    | 0    |
| SAMPLE1 q dcos [[/opt/F]SVxtc]anga/tcsds-1.2.36/lib64/libf]90f.so.1] [0]  | 0.14   | 0.14    | 3     | 0    |
| ISUMMARY] decomp 2d fft.fft 3d c2c [//vol0003/hp210193/data/benchmark3/version2decomp/FFT libraries/2decomp fft/src/fft fftw3.f90]      | 0.12   | 0.12    | 4     | 0    |
| SAMPLE] decomp 2d fft.fft 3d c2c [[/vol0003/hp210193/data/benchmark3/version2decomp/FFT libraries/2decomp fft/src/fft fftw3.f90] [540]] | 0.1    | 0.1     | 3     | 0    |
| SAMPLE] decomp 2d fft.fft 3d c2c [[/vol0003/hp210193/data/benchmark3/version2decomp/FFT libraries/2decomp fft/src/fft fftw3.f90] [498]  | 0.02   | 0.02    | 1     | 0    |
| SAMPLE] mmap [() {0]]   | 0.095  | 0.095   | 3     | 0    |
| SAMPLE] UNRESOLVED /opt/FJSVxos/mmm/lib64/libmpg.so.1   | 0.08   | 0.08    | 2     | 0    |
| [SUMMARY] decomp_2d.mem_merge_zy_complex_ [/vol0003/hp210193/data/benchmark3/version2decomp/FFT_libraries/2decomp_fft/src/./transpo     | 0.06   | 0.06    | 2     | 0    |
| [SUMMARY] decomp_2d.mem_merge_yx_real_ [{/vol0003/hp210193/data/benchmark3/version2decomp/FFT_libraries/2decomp_fft/src/./transpose_y   | 0.06   | 0.06    | 2     | 0    |
| SAMPLE]pthread_mutex_lock_full [{} {0}]   | 0.06   | 0.06    | 2     | 0    |
| [SUMMARY] decomp_2d.mem_split_xy_real_ [{/vol0003/hp210193/data/benchmark3/version2decomp/FFT_libraries/2decomp_fft/src/./transpose_x_t | 0.04   | 0.04    | 2     | 0    |
| [SUMMARY] decomp_2d.mem_split_yx_real_ [{/vol0003/hp210193/data/benchmark3/version2decomp/FFT_libraries/2decomp_fft/src/./transpose_y_t | 0.04   | 0.04    | 1     | 0    |
| [SUMMARY] decomp_2d.mem_merge_zy_real_ [{/vol0003/hp210193/data/benchmark3/version2decomp/FFT_libraries/2decomp_fft/src/./transpose_z   | 0.03   | 0.03    | 1     | 0    |
| [SUMMARY] decomp_2d.mem_split_yz_real_ [{/vol0003/hp210193/data/benchmark3/version2decomp/FFT_libraries/2decomp_fft/src/./transpose_y_t | 0.02   | 0.02    | 1     | 0    |
| SAMPLE] UNRESOLVED [vdso]   | 0.02   | 0.02    | 1     | 0    |
| [SAMPLE] qsort_r [{} {0}]   | 0.02   | 0.02    | 1     | 0    |
| [SAMPLE] UNRESOLVED /vol0003/hp210193/data/benchmark3/version2decomp/KleinGordon/2decomp_tau/sameer_final_run_384/Kg                    | 0.02   | 0.02    | 1     | 0    |
| [SAMPLE] msort_with_tmp.part.0 [{} {0}]   | 0.02   | 0.02    | 1     | 0    |
| ▼ ■ MPI_Alltoaliv()   | 13.685 | 13.685  | 138   | 0    |
| Television Context MPI_Alitoaliv()  | 0      | 13.657  | 457   | 0    |
| [SAMPLE] utofu_read_tcq_nocb [{/lib64/libtofucom.so} {0}]   | 3.37   | 3.37    | 113   | 0    |
| [SAMPLE] ompi_coll_mtofu_read_tcq_cq [{/opt/FJSVxtclanga/tcsds-1.2.36/lib64/libmpi.so.0.0.0} {0}]                                       | 3.078  | 3.078   | 101   | 0    |
| [SAMPLE] mca_btl_vader_component_progress [{/opt/FJSVxtclanga/tcsds-1.2.36/lib64/libmpi.so.0.0.0} {0}]                                  | 3.05   | 3.05    | 103   | 0    |
| [SAMPLE] ompi_coll_mtofu_alitoallv_doublespread [{/opt/FJSVxtclanga/tcsds-1.2.36/lib64/libmpi.so.0.0.0} {0}]                            | 1.37   | 1.37    | 48    | 0    |

# Figure: TAU's ParaProf shows the thread statistics table for MPI rank 0 for 2decomp&fft run on 384 ranks on Fugaku

Samar A. Aseeri

# TAU Profiling



Figure: TAU's ParaProf 3D visualization shows the shape of the overall profile of 2decomp&fft run on 384 ranks on Fugaku.

Samar A. Aseeri

Comparative Study of Profiling Tools on Fugaku Supercomputer

- Extrae-3.8.3 also collects performance information from binaries that are compiled with debug information.
- The instrumentation tool is configured to capture calls to the MPI library.
- The top MPI profile focus in the full execution reporting that the computation represents almost 70% of the execution time (Average value of Outside MPI) that seems a good percentage.
- With the MPI profile we can also measure that the balance of work between processes that in this case it is quite good (Avg/Max for Outside MPI is 93%) pointing that the problem is related with the communications (an analysis with Dimemas shows it can be related with data transfer).

# Extrae Profiling

| XX                                      |  |   |  |  |   |  |   |   |   |   |       |
|---|--|---|--|--|---|--|---|---|---|---|-------|
| i <b>c id</b> 3D                        | ् 🕿 🔳  | н Н Ш 2   | ¥ ∑ ½ ►  | Average 🗸 🕏  |   |  |   |   |   |   |       |
|   | Outside MPI  | MPI_Alltoall  | MPI_Finalize   | MPI_Comm_split   | MPI_Reduce  | MPI_Bcast  | MPI_Init  | MPI_Comm_size   | MPI_Comm_free   | MPI_Comm_rank   |       |
| Total                                   | 26,727.71 %  | 11,350.08 %   | 110.52 %   | 108.35 %   | 64.16 %   | 23.62 %  | 15.53 %   | 0.02 %  | 0.01 %  | 0.00 %  |       |
| Average                                 | 69.60 %  | 29.56 %   | 0.29 %   | 0.28 %   | 0.17 %  | 0.06 %   | 0.04 %  | 0.00 %  | 0.00 %  | 0.00 %  |       |
| Aaximum                                 | 71.03 %  | 31.69 %   | 0.29 %   | 0.31 %   | 0.32 %  | 0.09 %   | 0.04 %  | 0.00 %  | 0.00 %  | 0.00 %  |       |
| Minimum                                 | 67.33 %  | 28.23 %   | 0.08 %   | 0.18 %   | 0.04 %  | 0.00 %   | 0.00 %  | 0.00 %  | 0.00 %  | 0.00 %  |       |
| StDev                                   | 0.90 %   | 0.82 %  | 0.01 %   | 0.02 %   | 0.09 %  | 0.01 %   | 0.00 %  | 0.00 %  | 0.00 %  | 0.00 %  |       |
| Avg/Max                                 | 0.98   | 0.93  | 0.98   | 0.92   | 0.52  | 0.71   | 1.00  | 0.28  | 0.44  | 0.52  |       |
| <b>6</b> ×                              |  | Useful Dura   | tion @ ffte_384.   | prv  | ~ ^ 🔵   | <b>X</b> ×   |   |   |   |   |       |
|   |  |   |  |  |   | THREAD 1.11  |   |   |   |   |       |
|   | 0 us   | <mark></mark> 63  | ,751.51 - 63, <del>9</del> 96  | .64  | 26,186,735 us   | THREAD 1.38  | 57.1<br>84.1 0 us                                 |   |   |   | 26,11 |
| K ×<br>K = D 3D                         | Q 🛱 🔳  | ез<br>н н 11 р  | ,751.51 - 63, <del>99</del> 6<br>₩ Σ ½ ►   | .64<br>Average 🗸 🍫   | 26,186,735 93   | THREAD 1.38  | 4.1 0 us  |   |   |   | 26,18 |
| N ×                                     | Alitoali   | H H II 2<br>Outside MP1   | ,751.51 - 63,996<br>¥ ∑ ½ ⊾<br>MPI_Finalize  | .64<br>Average V &   | 26,186,735 us<br>MPI call profile<br>MPI_Reduce   | THREAD 1.38<br>THREAD 1.38<br>.c1 @ ffte_384<br>MPI_Bcast  | 57.1<br>84.1 () us<br>4.prv<br>MPI_Init           | MPI_Comm_size   | MPI_Comm_free   | MPI_Comm_rank   | 26,18 |
| C D 3D<br>Total                         | Q 😤 🔳<br>MPI_Alitoali<br>20,479.82 %                                   | 63<br>H H 11 7<br>Outside MP1<br>17,338.79 %  | ,751.51 - 63,996<br>Υ Σ ½ ►<br>MPI_Finalize<br>199.41 %  | .64<br>Average $\checkmark$ \$<br>MPI_Comm_split<br>195.50 %                                 | 26,186,735 93<br>MPI call profile<br>MPI_Reduce<br>115.77 %   | THREAD 1.38<br>THREAD 1.38<br>OF fite_384<br>MPI_Bcast<br>42.62 %  | 4.prv<br>MPI_Init<br>28.03 %                      | MPI_Comm_size   | MP1_Comm_free<br>0.01 %                               | MPI_Comm_rank<br>0.01 %                               | 26,18 |
| K D 3D<br>Total<br>Average              | C. C                               |   | ,751.51 - 63,996<br>★ Σ ½ ►<br>MPI_Finalize<br>199.41%<br>0.52%  | .44 Average V & MPI_Comm_split 195.50 % 0.51 %   | 26,186,733 un<br>MPI call profile<br>MPI_Reduce<br>115.77 %<br>0.30 %                               | HREAD 1.33<br>THREAD 1.33<br>MPI_Bcast<br>42.62 %<br>0.11 %  | 4.prv<br>MPI_Init<br>28.03 %<br>0.07 %            | MPI_Comm_size<br>0.04 %<br>0.00 %                     | MPI_Comm_free<br>0.01 %<br>0.00 %                     | MPI_Comm_rank<br>0.01 %<br>0.00 %                     | 26,18 |
| K Z D 3D<br>Total<br>Average<br>Maximum | MPI_Alitoali<br>20,479.82 %<br>53.33 %<br>57.18 %                      | H H III 2<br>Outside MP1<br>17,338.79 %<br>45.15 %<br>47.72 %   | ,753.53 - 63,990<br>★ Σ ½ ►<br>MPI_Finalize<br>199.41 %<br>0.52 %<br>0.53 %  | Average V &<br>MPI_Comm_split<br>195.50 %<br>0.51 %<br>0.55 %                                | 26,186,733 us<br>MPI call profile<br>MPI_Reduce<br>115.77 %<br>0.30 %<br>0.58 %                     | HREAD 1.25<br>THREAD 1.33<br>MPI_Bcast<br>42.62 %<br>0.11 %<br>0.16 %                                      | 4.prv<br>MPI_Init<br>28.03 %<br>0.07 %<br>0.07 %  | MPI_Comm_size<br>0.04 %<br>0.00 %<br>0.00 %           | MPI_Comm_free<br>0.01 %<br>0.00 %                     | MPI_Comm_rank<br>0.01 %<br>0.00 %<br>0.00 %           | 26,18 |
| Total<br>Average<br>Maximum<br>Minimum  | MPI_Alitoall<br>20,479.82 %<br>53.33 %<br>57.18 %<br>50.94 %           | 63<br>H H III 2<br>Outside MPI<br>17,338.79 %<br>45.15 %<br>47.72 %<br>41.05 %  | ,751.51 - 63,990   | Average V &  | 26,186,733 sr<br>MPL call profile<br>MPLReduce<br>115.77 %<br>0.30 %<br>0.58 %<br>0.08 %            | HREAD 1.25<br>THREAD 1.35<br>.c1 @ ffte_384<br>MPI_Bcast<br>42.62 %<br>0.11 %<br>0.16 %<br>0.00 %          | MPI_Init<br>28.03 %<br>0.07 %<br>0.07 %<br>0.00 % | MPI_Comm_size<br>0.04 %<br>0.00 %<br>0.00 %           | MPI_Comm_free<br>0.01 %<br>0.00 %<br>0.00 %           | MPI_Comm_rank<br>0.01 %<br>0.00 %<br>0.00 %           | 26,18 |
| Total<br>Average<br>Maximum<br>StDev    | MPI_Alitoall<br>20,479.82 %<br>53.33 %<br>57.18 %<br>50.94 %<br>1.48 % | H         H         III         2           Outside MPI         17,338.79 %         45.15 %         47.72 %           41.05 %         1.63 %         1.63 % | x    x     x     x     x | Average V<br>AVerage V<br>MPI_Comm_split<br>195.50 %<br>0.51 %<br>0.52 %<br>0.32 %<br>0.04 % | 26,186,733 sr<br>MPI call profile<br>MPI_Reduce<br>115.77 %<br>0.30 %<br>0.58 %<br>0.08 %<br>0.08 % | HREAD 1.25<br>THREAD 1.35<br>C1 @ ffte_384<br>MPI_Bcast<br>42.62 %<br>0.11 %<br>0.16 %<br>0.00 %<br>0.03 % | MPI_Init<br>28.03 %<br>0.07 %<br>0.00 %           | MPI_Comm_size<br>0.04 %<br>0.00 %<br>0.00 %<br>0.00 % | MPJ_Comm_free<br>0.01 %<br>0.00 %<br>0.00 %<br>0.00 % | MPI_Comm_rank<br>0.01 %<br>0.00 %<br>0.00 %<br>0.00 % | 26,18 |

Figure: Extrae MPI instrumentation of FFTE version on Fugaku

Samar A. Aseeri

- Focusing on the goal to compare the executions with 384 cores, we report a quick analysis that can be done with Paraver for the performance data collected by Extrae.
- As the two versions execute a different number of iterations, we focused the analysis on metrics for the full run.
- Both timelines are on the same scale reflecting that FFTE is much faster than 2Decomp&FFT.

# Analysis using Extrae to Compare Executions with 384 tasks



Figure: FFTE has shorter runtime than 2Decomp&FFT, but has less synchronized MPI calls.

- FIPP and FAPP have higher overhead.
- Unlike TAU and Extrae, FIPP and FAPP do not support dynamic instrumentation of a binary executable with debug information, the instrumentation for FIPP and FAPP must be done during compilation.
- The source codes need to be annotated with the commands fipp\_start/fapp\_start and ended with the fipp\_stop/fapp\_stop.
- There is no special tool to visualize FIPP data on Fugaku.
- FIPP produces a summary file that enables one to determine processor imbalance in the annotated region.

# **FIPP** Profiling

| Time | statistics         |               |           |           |    |
|------|--------------------|---------------|-----------|-----------|----|
|      | Elapsed(s)         | User(s)       | System(s) |           |    |
|      | 38.5575            |               |           | Applicati | on |
|      | 38.5575            |               |           | Process   | 5  |
|      | 38.5567            |               |           | Process   | 7  |
|      | 38.5563            |               |           | Process   | 4  |
|      | 38.5561            |               |           | Process   | 3  |
|      | 38.5557            |               |           | Process   | 0  |
|      | 38.5540            |               |           | Process   | 2  |
|      | 38.5540            |               |           | Process   | 1  |
|      | 38.5419            |               |           | Process   | 6  |
|      | 38.5378            |               |           | Process   | 8  |
|      | 38.5216            |               |           | Process   | 12 |
|      | 38.5183            |               |           | Process   | 11 |
|      | 38.5058            |               |           | Process   | 13 |
|      | 38.5055            |               |           | Process   | 9  |
|      | 38.4910            |               |           | Process   | 10 |
|      | 38.4784            |               |           | Process   | 14 |
|      | 38.4605            |               |           | Process   | 16 |
|      |                    |               |           |           |    |
| erfo | rmance monitor eve | nt:statistics |           |           |    |

Figure: Head of FIPP Profiling data for 2decomp&fft

Samar A. Aseeri

# **FAPP** Profiling



Figure: FAPP CPU Performance Analysis Report for the 2Decomp&FFT executed on 384 tasks for reading the Pa1,Pa2,Pa3,Pa4 and Pa5 reports to obtain the Brief Report

Samar A. Aseeri

Comparative Study of Profiling Tools on Fugaku Supercomputer

- 2 Background & Related Work
- 3 Profiling Results: A Comprehensive Overview
- 4 Evaluating the Results of Profiling
- Summary & Conclusions

#### 6 Further Work

- We have run three test cases of 284<sup>3</sup> 768<sup>3</sup> and 1536<sup>3</sup> grid sizes, with GNU and with Fujitsu compilers but Fujitsu was superior as the three run average gave 12.56851788998271 while the GNU average gave 38.865136866666641 for the 7683 case.
- Consequently, for our experiments we use default Fugaku environment "lang/tcsds- 1.2.36(default)" to compile our MPI code.
- For this, we sought to consider three metrics: run-time overhead, memory usage, size of generated data.

# Tool Overhead Comparisons

#### Strong Scaling of Klein Gordon + small-torus mode Error bars show standard deviation in runtime



Figure: Performance Tools Comparison with Small-Torus mode

Comparative Study of Profiling Tools on Fugaku Supercomputer

30 / 42

# Overhead Memory Measurements



2DECOMP&FFT

Figure: Performance Tools memory-check comparison of all cases for 2Decomp&FFT Library

# Overhead Memory Measurements



Figure: Performance Tools memory-check comparison of all cases for FFTE Library

- Size of data generated from running the executable with the examined profiling tools can be found in Tables in the upcoming slides.
- TAU supports compression and normalization of profile files using the paraprof -pack command to generate files in TAU's packed profile format.
- The ppk format is very efficient in reducing the size of the profile dataset and can result in several orders of magnitude reduction in data sizes.

Table: Table comparing data size generated for parallel FFT for a  $768^3$  problem size (30 forward and backward FFTs) using 2Decomp&FFT library executed with the Extrae, TAU and FIPP tools on Fugaku.

| Nodes | Cores | 2Decomp&FFT | 2Decomp&FFT | 2Decomp&FFT | 2Decomp&FFT | 2Decomp&FFT |
|-------|-------|-------------|-------------|-------------|-------------|-------------|
|       |       |             | Extrae      | TAU         | FIPP        | FAPP        |
| 8     | 384   | 623.23K     | 79.18M      | 11.29M      | 9.54M       | 23.00M      |
| 16    | 768   | 629.42K     | 281.12M     | 22.80M      | 17.42M      | 11.80M      |
| 32    | 1536  | 642.29K     | 730.95M     | 37.13M      | 28.94M      | 6.20M       |
| 64    | 3072  | 668.33K     | 1.42G       | 64.37M      | 55.17M      | 3.40M       |
| 128   | 6144  | 720.44K     | 2.82G       | 115.35M     | 99.73M      | 2.00M       |
| 256   | 12288 | 826.97K     | 5.59G       | 219.98M     | 197.52M     | 1.30M       |

Table: Table comparing data size generated for parallel FFT for a 768<sup>3</sup> problem size (30 forward and backward FFTs) using FFTE library with executed with the Extrae, TAU and FIPP tools on Fugaku.

| Nodes | Cores | FFTE    | FFTE + Extrae | FFTE + TAU | FFTE + FIPP | FFTE + FAPP |
|-------|-------|---------|---------------|------------|-------------|-------------|
| 8     | 384   | 385.13K | 20.91M        | 20.31M     | 3.11M       | 17.15M      |
| 16    | 768   | 391.28K | 42.11M        | 30.32M     | 5.12M       | 8.76M       |
| 32    | 1536  | 404.13K | 91.72M        | 42.91M     | 8.54M       | 4.56M       |
| 64    | 3072  | 430.24K | 188.33M       | 65.09M     | 15.20M      | 2.47M       |
| 128   | 6144  | 482.45K | 439.29M       | 101.76M    | 28.28M      | 1.42M       |
| 256   | 12288 | 589.08K | 920.25M       | 172.14M    | 57.52M      | 916.75K     |

Table: Table comparing data sizes for TAU's packed profile format (ppk) for 2Decomp&FFT and FFTE library generated for parallel FFT for a 768<sup>3</sup> problem size (30 forward and backward FFTs) on Fugaku.

| Nodes | Cores | 2Decomp&FFT + TAU (ppk) | FFTE + TAU (ppk) |
|-------|-------|-------------------------|------------------|
| 8     | 384   | 461K                    | 574K             |
| 16    | 768   | 840K                    | 862K             |
| 32    | 1536  | 1.3M                    | 1.2M             |
| 64    | 3072  | 2.3M                    | 1.8M             |
| 128   | 6144  | 3.9M                    | 2.9M             |
| 256   | 12288 | 7.0M                    | 4.5M             |

- 2 Background & Related Work
- Operation of the second sec
- Evaluating the Results of Profiling
- 5 Summary & Conclusions

#### 6 Further Work

# Summary & Conclusions

- This study examines the overhead of performance tool measurements when profiling Fast Fourier Transform (FFT) based solvers for the Klein Gordon equation. The tools examined are TAU and Extrae in addition to the Fujitsu provided Fugaku profiling tools FAPP and FIPP.
- Performance measurement is an important part of improving the efficiency of parallel programs for a specific architecture.
- Fugaku's small torus mode is shown to significantly improve performance of communication intensive codes by ensuring communication from other jobs is not intermingled with communication of the job.
- This paper examines the performance of programs that heavily utilize the FFT, a communication intensive algorithm, by minimizing communication interference from other running jobs.

- It shows that profiling and tracing using binary instrumentation can capture many aspects of a program, but knowledge of the program structure is still required to interpret the resulting performance measurements.
- The paper focuses on performance measurement using four tools and two FFT libraries, but there are other tools and FFT libraries that can also be used.
- The paper also presents memory overhead measurements and profiling result storage requirements.
- The Klein Gordon solver is used as a mini application benchmark for a parallel program that solves a partial differential equation.

- 2 Background & Related Work
- 3 Profiling Results: A Comprehensive Overview
- Evaluating the Results of Profiling
- Summary & Conclusions



- When McKernel is available on Fugaku, it is important to measure the effect of OS noise reduction on performance of bulk synchronous programs.
- Further work involves profiling other parallel FFT libraries and additional profiling and tracing tools, examining other exascale hardware architectures, and examining other numerical methods to aid in performance prediction and optimal matching of algorithms to hardware architecture.
- Finally, GPU-based FFT libraries should be analyzed for their potential to improve performance on exascale systems.

# Acknowledgments



